

LA-UR-20-20324

Approved for public release; distribution is unlimited.

Title: Robust Bayesian change detection for cyber-security applications

Author(s): Hallgren, Karl Lars Yvon
Turcotte, Melissa

Intended for: Student presentation

Issued: 2020-01-13

Disclaimer:

Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by Triad National Security, LLC for the National Nuclear Security Administration of U.S. Department of Energy under contract 89233218CNA000001. By approving this article, the publisher recognizes that the U.S. Government retains nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

Robust Bayesian change detection for cyber-security applications

Karl Hallgren

Mentor at LANL: Melissa Turcotte[†]

PhD supervisors: Nicholas Heard* and Niall Adams*

[†] Advanced Research in Cyber-Systems (A-4), Los Alamos National Laboratory

* Department of Mathematics, Imperial College London (UK)

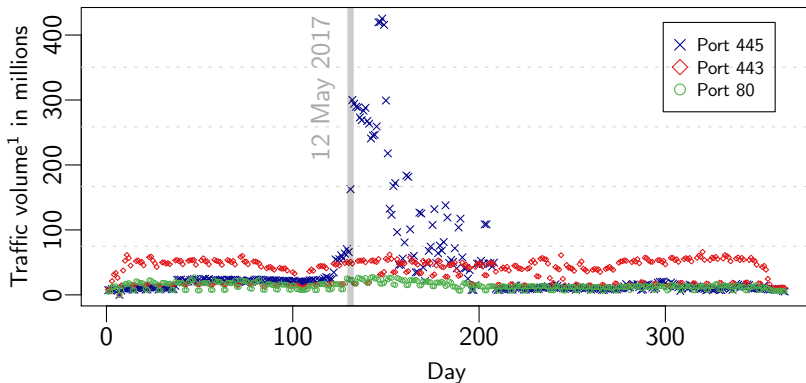
August, 2019



Motivation for change detection

A cyber attack typically changes the behaviour of the target network.

For example, by exploiting a Windows vulnerability the 2017 WannaCry ransomware attack led to spikes of activity on port 445.



¹recorded by a network router at Imperial College London

Motivation for change detection

To detect the presence of a network intrusion, it can be informative to monitor for changes in the high-volume data sources which are collected inside an enterprise computer network

- NetFlow - summaries of connection between devices
- Authentication events
- Host event logs

Large data set derived from the operational network environment at LANL

- Unified Host and Network Data Set (Turcotte et al., 2017)
<https://csr.lanl.gov/data/2017.html>

Bayesian changepoint analysis

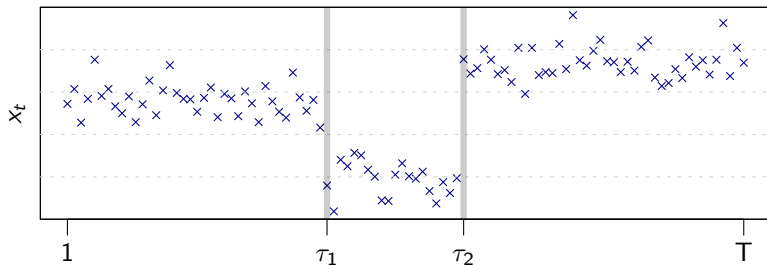
An unknown number k of changepoints, whose positions are denoted by $\tau_{1:k} = (\tau_1, \dots, \tau_k)$, split the data x_1, \dots, x_T into $k + 1$ segments such that within each segment j

$$x_t \stackrel{\text{iid}}{\sim} f(\cdot | \theta_j) \quad \forall t \in \{\tau_{j-1}, \dots, \tau_j - 1\}$$

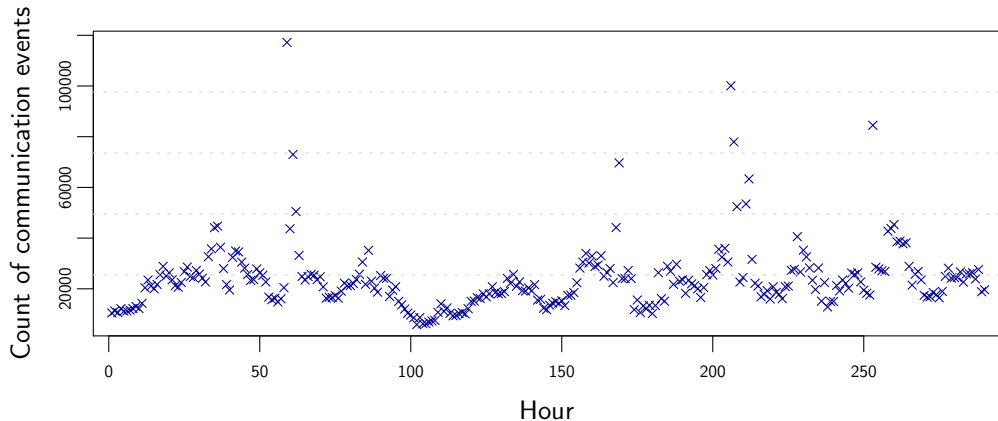
for some segment parameter θ_j .

A priori each time point is assumed to be a changepoint with probability p so that

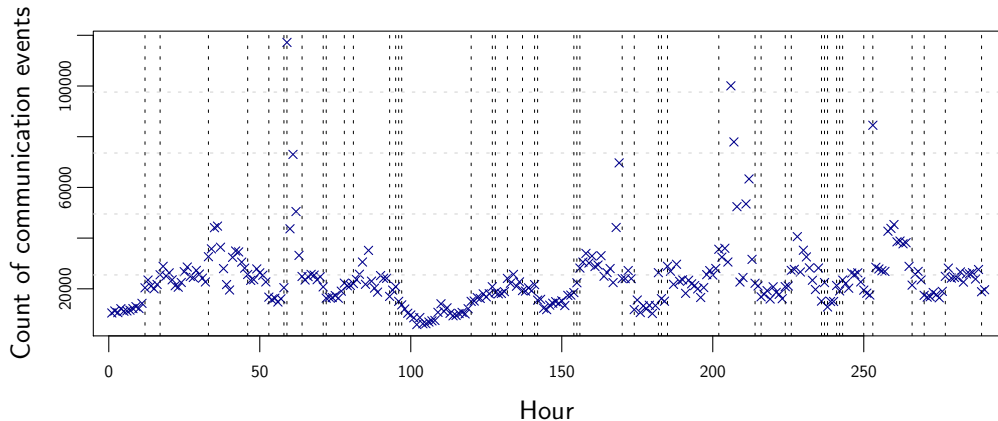
$$\pi(k, \tau_{1:k}) = p^k (1 - p)^{T-1-k} \propto \{p/(1 - p)\}^k$$



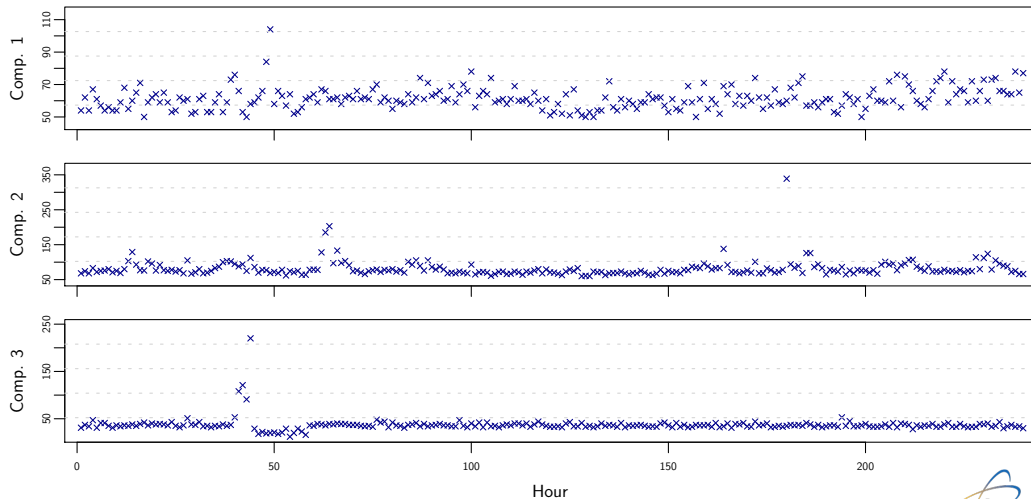
Difficulties applying traditional change detection methods



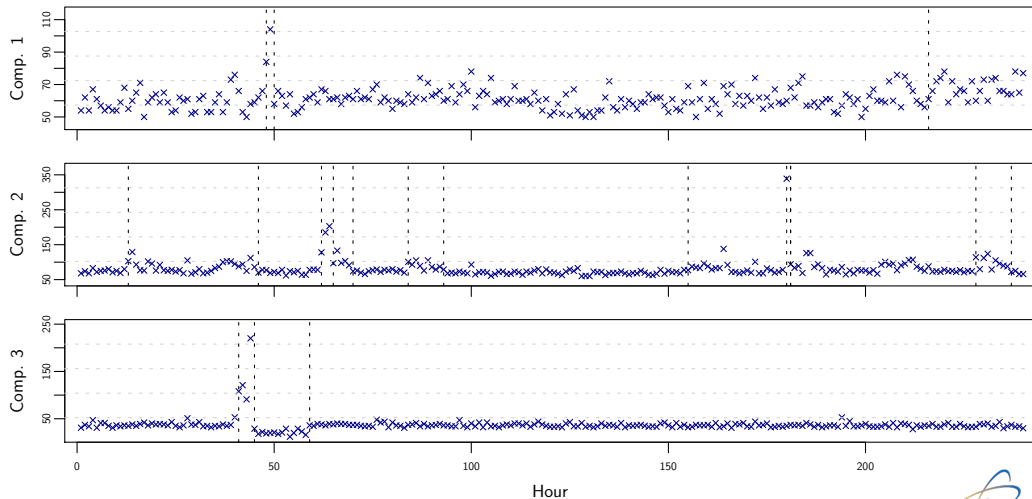
Difficulties applying traditional change detection methods



Difficulties applying traditional change detection methods



Difficulties applying traditional change detection methods



Difficulties applying traditional change detection methods

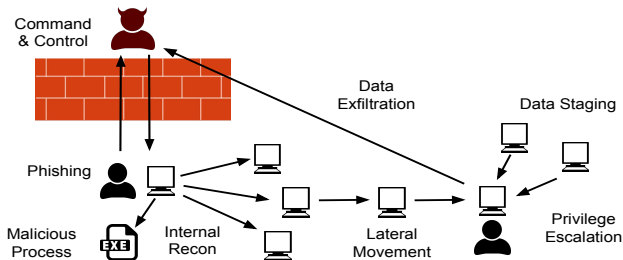
Most traditional changepoint detection methods will fail to characterise what cyber-security analysts mean by a change and consequently fit many more changepoints than preferable.

1. Classical changepoint models are not robust to seasonal variations, gradual drifts and other *normal* dynamic phenomena one may observe in cyber data.
2. Suspicious changes in the behaviour of a computer within a network can be observed under non-attack conditions.

Robust Bayesian change detection for cyber-security applications

Therefore it is relevant to develop methods to reduce false positive detections in changepoint analysis for cyber-security applications.

1. Building changepoint models which are robust to *normal* dynamic phenomena one may observe in cyber data, yet still tractable.
2. Combining evidence from multiple sources to identify patterns of changes which are a priori likely to correspond to the kill chain of an attack (Hutchins et al., 2011; Sexton et al., 2015).



Part 1. Building robust segment models

Relaxing the assumption of exchangeability

For each segment, assuming the data are iid conditional on some segment-specific parameter is equivalent to assuming that the data are infinitely exchangeable by De Finetti's representation theorem (Bernardo and Smith, 1993).

A potentially infinite sequence x_1, x_2, \dots is said to be *infinitely exchangeable* under a probability measure F if, for every n and any set of indices $I_n = \{i_1, \dots, i_n\}$, we have

$$F(x_{i_1}, \dots, x_{i_n}) = F(x_{\sigma(i_1)}, \dots, x_{\sigma(i_n)}),$$

for any permutation σ on I_n .

→ This assumption of complete symmetry within segments is too restrictive for cyber-security applications.

Segment models for non-exchangeable data

It is of interest to build segment models which may admit exchangeability but also other weak forms of dependence, so that the corresponding changepoint models are more flexible (yet still tractable).

Moving-sum changepoint model

Within a generic segment, suppose the observed data x_1, \dots, x_n satisfy

$$x_t = \sum_{i=0}^m y_{t-i},$$

where (y_t) are latent random variables such that

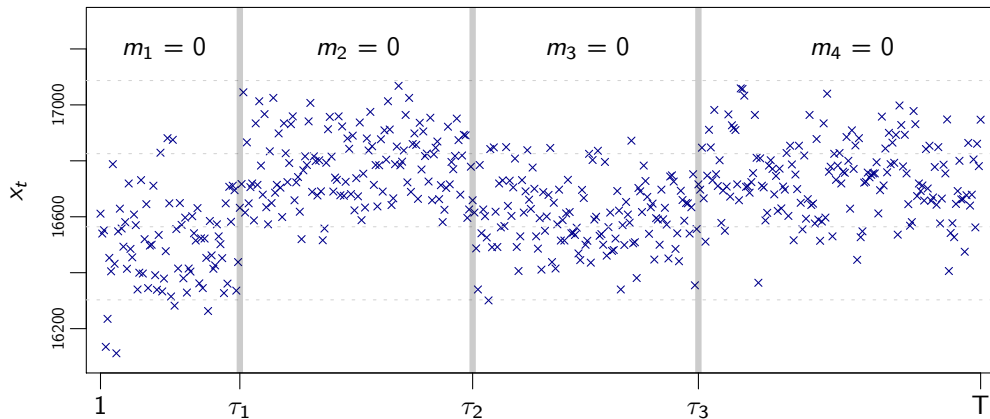
$$y_t \stackrel{\text{iid}}{\sim} f(\cdot | \theta, m)$$

for some *unknown* segment-specific parameters $\theta \in \Theta$ and $m \geq 0$.

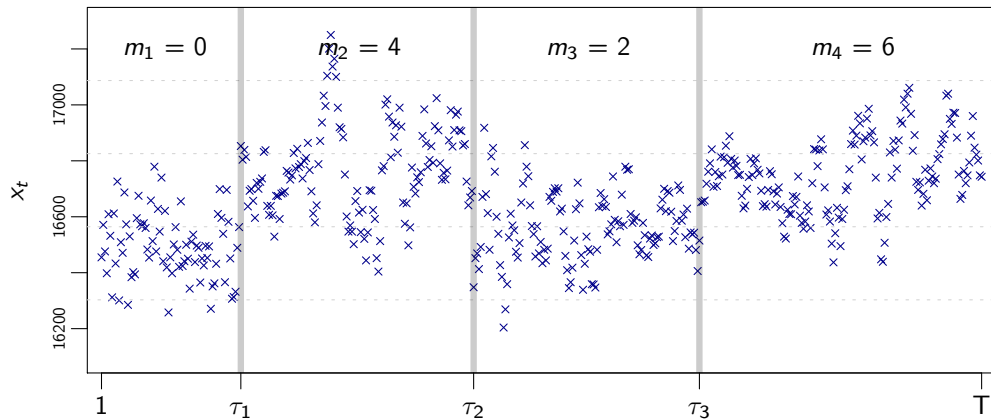
This approach allows m -dependence within segments, whilst maintaining a desired marginal distribution in the class of convolution-closed infinitely divisible distributions.

- E.g., if $y_t \stackrel{\text{iid}}{\sim} \text{Poisson}(\theta/(m+1))$ then $x_t \sim \text{Poisson}(\theta)$ with m -dependence.

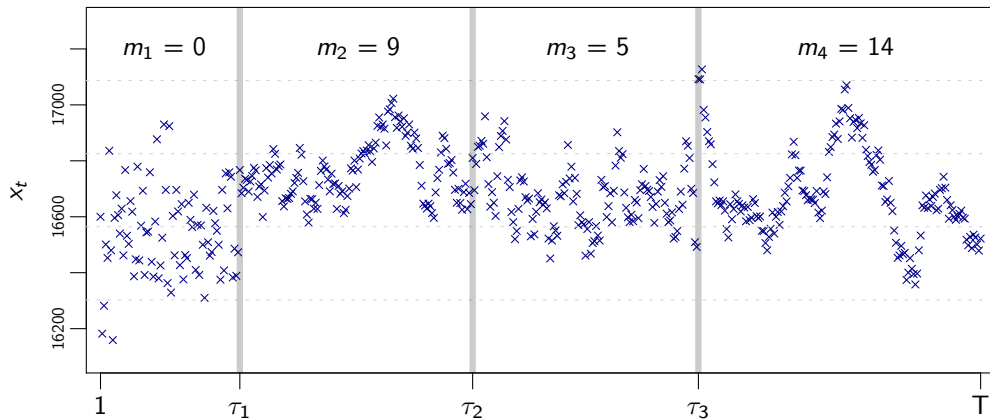
Simulations from the moving-sum changepoint model



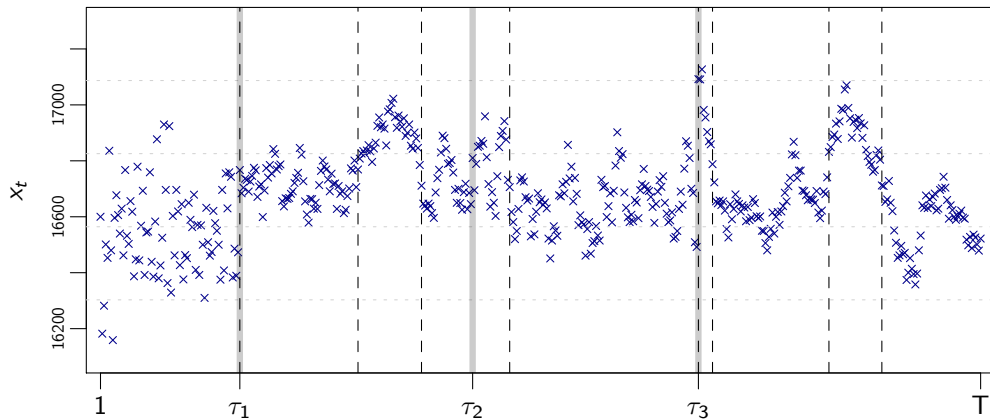
Simulations from the moving-sum changepoint model



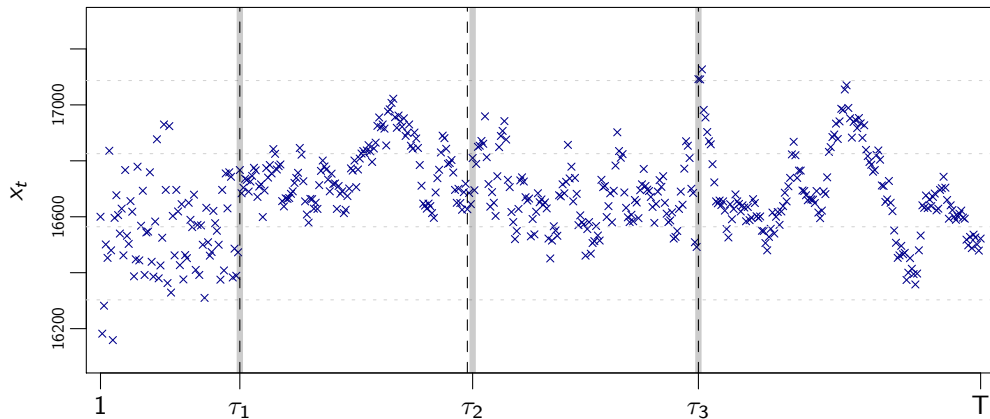
Simulations from the moving-sum changepoint model



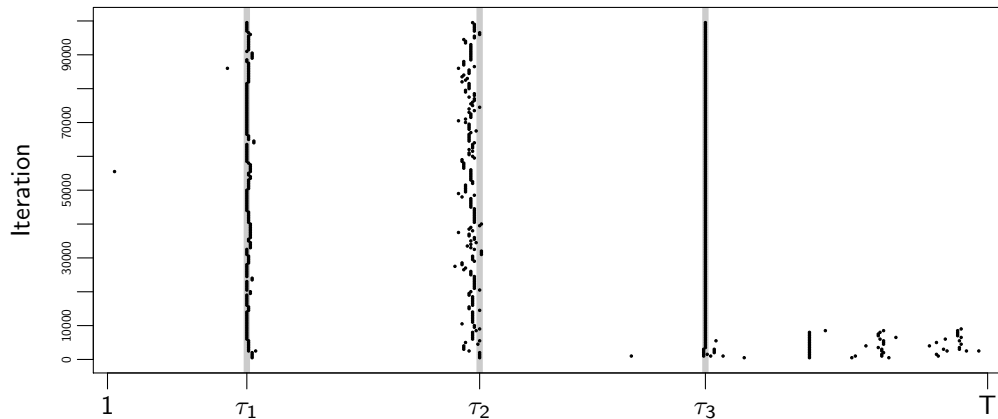
Simulations from the moving-sum changepoint model



Simulations from the moving-sum changepoint model



Simulations from the moving-sum changepoint model



Part 2. Towards kill chain detection

Inadequate prior distribution for the changepoints

The aim is to perform changepoint analysis on L time series of length T , denoted by $\mathbf{y} = (y_{\ell,t})$.

Let $\mathbf{s} = (s_{\ell,t})$ be a matrix indicating the positions of the changepoints,

$$s_{\ell,t} = \begin{cases} 1 & \text{if } t \text{ is a changepoint for the } \ell\text{-th time series} \\ 0 & \text{otherwise.} \end{cases}$$

A priori assumed that

$$\pi(\mathbf{s}) = \prod_{\ell,t} \{p^{s_{\ell,t}}(1-p)^{1-s_{\ell,t}}\} \propto \prod_{\ell,t} \left\{ \frac{p}{1-p} \right\}^{s_{\ell,t}}$$

One issue is that $\pi(\mathbf{s}) = \pi(\mathbf{s}')$, where

$$\mathbf{s} = \begin{pmatrix} 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \end{pmatrix} \quad \mathbf{s}' = \begin{pmatrix} 0 & 0 & 0 & 0 & \mathbf{1} & 0 & \mathbf{1} & 0 \\ 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Encoding prior knowledge on cyber-attack propagation

Cyber-analysts can specify prior knowledge on cyber-attack propagation by defining the neighbourhoods of the indices ℓ and t , denoted by $\partial(\ell)$ and $\partial(t)$, respectively.

For example, if $\ell, \ell_1, \ell_2, \dots, \ell_L$ denote the indices of time-series representing the process activity of computers which are known to be *peers* on the network, we may have

$$\{\ell_1, \dots, \ell_L\} \subset \partial(\ell) \quad \text{and} \quad \partial(t) = \{t-1, t, t+1\}.$$

To encode the structure in the neighbourhood of (ℓ, t) , define

$$w_{\ell,t} \equiv w_{\ell,t}(\mathbf{s}) = \left(1 + \sum_{\ell' \in \partial(\ell)} \mathbb{1}\{1 \in s_{\ell', \partial(t)}\} \right) / (1 + |\partial(\ell)|)$$

For example,

$$w_{\ell,t} = 1/3 \quad \text{if} \quad \mathbf{s} = \begin{pmatrix} 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \textcircled{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 \end{pmatrix}$$

Encoding prior knowledge on cyber-attack propagation

Cyber-analysts can specify prior knowledge on cyber-attack propagation by defining the neighbourhoods of the indices ℓ and t , denoted by $\partial(\ell)$ and $\partial(t)$, respectively.

For example, if $\ell, \ell_1, \ell_2, \dots, \ell_L$ denote the indices of time-series representing the process activity of computers which are known to be *peers* on the network, we may have

$$\{\ell_1, \dots, \ell_L\} \subset \partial(\ell) \quad \text{and} \quad \partial(t) = \{t-1, t, t+1\}.$$

To encode the structure in the neighbourhood of (ℓ, t) , define

$$w_{\ell,t} \equiv w_{\ell,t}(\mathbf{s}) = \left(1 + \sum_{\ell' \in \partial(\ell)} \mathbb{1}\{1 \in s_{\ell', \partial(t)}\} \right) / (1 + |\partial(\ell)|)$$

For example,

$$w_{\ell,t} = 3/3 \quad \text{if} \quad \mathbf{s} = \begin{pmatrix} 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \textcircled{1} & 0 & 0 & 0 & 0 \end{pmatrix}$$

Encoding prior knowledge on cyber-attack propagation

Cyber-analysts can specify prior knowledge on cyber-attack propagation by defining the neighbourhoods of the indices ℓ and t , denoted by $\partial(\ell)$ and $\partial(t)$, respectively.

For example, if $\ell, \ell_1, \ell_2, \dots, \ell_L$ denote the indices of time-series representing the process activity of computers which are known to be *peers* on the network, we may have

$$\{\ell_1, \dots, \ell_L\} \subset \partial(\ell) \quad \text{and} \quad \partial(t) = \{t-1, t, t+1\}.$$

To encode the structure in the neighbourhood of (ℓ, t) , define

$$w_{\ell,t} \equiv w_{\ell,t}(\mathbf{s}) = \left(1 + \sum_{\ell' \in \partial(\ell)} \mathbb{1}\{1 \in s_{\ell', \partial(t)}\} \right) / (1 + |\partial(\ell)|)$$

For example,

$$w_{\ell,t} = 3/3 \quad \text{if} \quad \mathbf{s} = \begin{pmatrix} 0 & 0 & \mathbf{1} & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \textcircled{1} & 0 & 0 & 0 & 0 \end{pmatrix}$$

Informative prior distribution for the changepoints

Consider a prior distribution for \mathbf{s} which takes into account a priori knowledge on cyber-attack propagation,

$$\pi(\mathbf{s}) \propto \prod_{\ell, t} \left\{ \left(\frac{p}{1-p} \right) \phi_{\lambda}(w_{\ell, t}) \right\}^{s_{\ell, t}}$$

for some non-decreasing function $\phi_{\lambda} : (0, 1] \rightarrow (0, 1]$ such that $\phi_{\lambda}(1) = 1$.

Now,

$$\pi(\mathbf{s})/\pi(\mathbf{s}') = \frac{\phi_{\lambda}^3(1)}{\phi_{\lambda}^2(2/3)\phi_{\lambda}(1/3)} \geq 0,$$

with

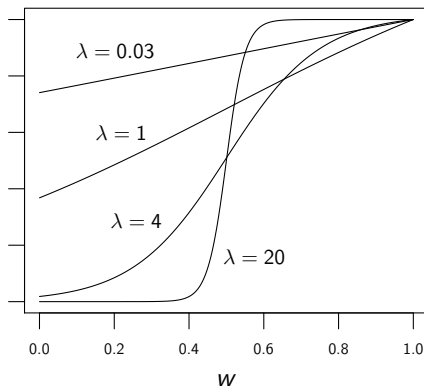
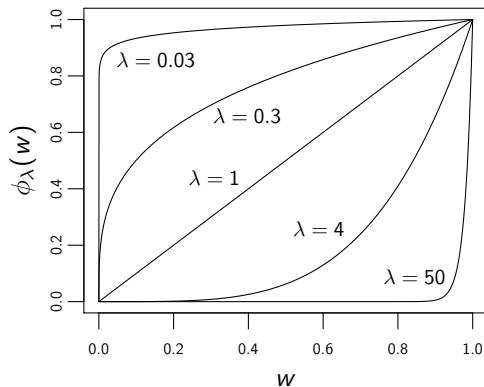
$$\mathbf{s} = \begin{pmatrix} 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \end{pmatrix} \quad \mathbf{s}' = \begin{pmatrix} 0 & 0 & \mathbf{1} & 0 & 0 & 0 & \mathbf{1} & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Flexible informative prior distribution for the changepoints

Possible choices for the link function ϕ_λ include:

$$\phi_\lambda(w) = w^\lambda$$

$$\phi_\lambda(w) = \frac{1 + \exp(-\lambda)}{1 + \exp(-\lambda[2w - 1])}$$



Let \mathbf{s} be binary matrix with $s_{\ell,t} = 0$, and let \mathbf{s}' be identical to \mathbf{s} but with $s_{\ell,t} = 1$.

Gibbs sampling is possible since the full conditional distribution of $s_{\ell,t}$ is available,

$$s_{\ell,t} \sim \text{Bernoulli} \left(\frac{Q}{1 + Q} \right),$$

where

$$Q = \frac{f(\mathbf{y} | \mathbf{s}')}{f(\mathbf{y} | \mathbf{s})} \frac{\pi(\mathbf{s}')}{\pi(\mathbf{s})},$$

but poor mixing and convergence issues are to be expected since the parameters will typically be highly correlated.

Conclusion

Research presented in this talk was supported by the Laboratory Directed Research and Development program of Los Alamos National Laboratory (New Mexico, USA) under project number 20180607ECR.

- Bernardo, J. M. and Smith, A. F. M. (1993). *Bayesian theory*. Wiley series in probability and mathematical statistics. John Wiley, New York; Chichester.
- Hutchins, E., Cloppert, M., and Amin, R. (2011). Intelligence-driven computer network defense informed by analysis of adversary campaigns and intrusion kill chains. *Leading Issues in Information Warfare and Security Research*, 1.
- Sexton, J., Storlie, C. B., and Neil, J. (2015). Attack chain detection. *Statistical Analysis and Data Mining*, 8:353–363.
- Turcotte, M. J. M., Kent, A. D., and Hash, C. (2017). Unified Host and Network Data Set. *ArXiv e-prints*. 1708.07518.

Alternative representation of the moving-sum segment model

Recall $x_t = \sum_{i=0}^m y_{t-i}$ for all t . Equivalently,

$$y_t = y_{t-(m+1)} + x_t - x_{t-1}, \quad \text{for } t = 2, \dots, n, \quad (1)$$

and $y_1 = x_1 - (y_{-m+1} + \dots + y_0)$.

The first m latent variables $\gamma_{1:m} = (\gamma_1, \dots, \gamma_m)$, with

$$\gamma_r = y_{-m+r}, \quad \text{for } r = 1, \dots, m,$$

may be seen as the unknown initial conditions of a stochastic difference equation defined by (1).

Conditional on the $\gamma_{1:m}$, there is a one-to-one deterministic transformation, denoted by Υ , between $x_{1:n}$ and $y_{1:n}$ with unit Jacobian,

$$y_{1:n} = \Upsilon(x_{1:n} \mid \gamma_{1:m}),$$

which can be obtained explicitly by iterating (1).

Conditional likelihood of the observed data

As a result, for each segment resulting from the changepoints, the conditional likelihood of the data

$$f(x_{1:n} | \theta, m, \gamma_{1:m}) = f(y_{1:n} | \theta, m, \gamma_{1:m}) \times 1.$$

is tractable, and therefore if we treat $(\theta, m, \gamma_{1:m})$ as unknown segment-specific parameters with prior density

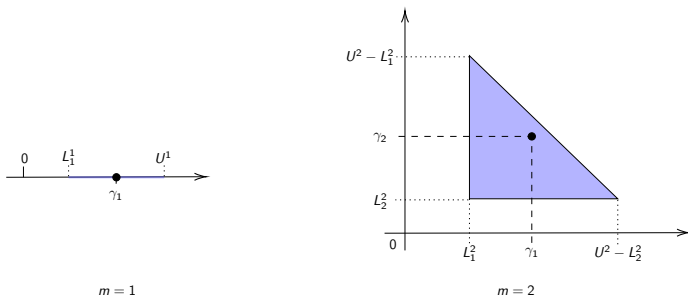
$$\pi(\theta)\pi(m) \prod_{r=1}^m f(\gamma_r | \theta, m),$$

then an expression for the joint posterior distribution of the changepoints and the segment parameters is immediately available via Bayes' theorem.

The segment parameter space

Let $\mathcal{Y}_m \equiv \mathcal{Y}_m(x_{1:n})$ denote the set of sequences $\gamma_{1:m}$ such that y_t belongs to the support of $f(y | m, \theta)$, for all $t = 1, \dots, n$.

\mathcal{Y}_m depends on the observed data when the support of $f(y | m, \theta)$ is bounded.



The structure of the constrained parameter space can be exploited to design a sampling strategy.